

# A Block Coordinated Update Method for Beam-Based Massive MIMO Downlink Scheduling Based on Statistical CSI

Zhiyuan Jiang, Sheng Zhou, and Zhisheng Niu, *Fellow, IEEE*  
 Tsinghua National Laboratory for Information Science and Technology,  
 Department of Electronic Engineering, Tsinghua University, Beijing, 100084, China  
 {zhiyuan, sheng.zhou, niuzhs}@tsinghua.edu.cn

**Abstract**—In this paper, the joint user and beam scheduling problem in beam-based massive multiple-input multiple-output (MIMO) systems is formulated based on the Lyapunov-drift optimization framework and the optimal scheduling policy is given in a closed-form. To address the weighted sum rate maximization problem (mixed integer programming) arisen in the Lyapunov-drift maximization, an algorithm based on the block coordinated update is proposed and proved to converge to the global optimum of the relaxed convex problem. In order to make the scheduling decisions based only upon statistical channel state information (CSI), asymptotic expressions of the downlink broadcast channel capacity are derived. Simulation results based on widely-adopted spatial channel models are given, which show that the proposed scheme is close to the optimal scheduling scheme, and outperforms the state-of-the-art beam selection schemes.<sup>1</sup>

## I. INTRODUCTION

Under the assumption that full digital signal processing for each antenna is performed at the base station (BS) side in massive multiple-input multiple-output (MIMO) based wireless communication systems, the system performance has been extensively investigated, e.g., in [1], [2]. However, it is widely accepted that full digital signal processing implementation encounters very severe challenges in practice, on account of the following impediments. The radio-frequency (RF) chain hardware cost, including e.g., low-noise amplifier, analog-digital-converter (ADC), power amplifier and etc., and corresponding power consumptions are high [3]. Baseband spatial signal processing complexity is extremely increased. Moreover, there are some other practical considerations which are system specific, e.g., fronthaul interface capacity limitation in cloud radio access networks (C-RAN), and the channel state information (CSI) acquisition overhead in frequency-division-duplexing (FDD) systems [4].

In view of these challenges, architectures with low RF- and processing-complexity have been proposed, see e.g., in [3], [5]–[7]. The existing literature can be divided into three categories. The first is *hybrid beamforming*, which adopts an RF front end with an analog beamforming module such that the number of RF chains is significantly reduced. The recently proposed *beam-space MIMO* architecture [6] adopts lens

<sup>1</sup>This work is sponsored in part by the Nature Science Foundation of China (No. 61701275, No. 91638204, No. 61571265, No. 61621091), the China Postdoctoral Science Foundation, and Hitachi R&D Headquarter..

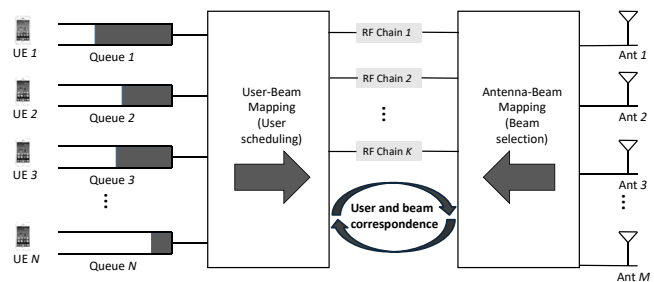


Fig. 1. An illustration of the beam-based massive MIMO systems where the user scheduling and beam selection are correlated.

antenna arrays which acts analogously like a lens focusing on light beams from different directions, and effectively as a directional analog beamforming module. It transforms the signal into the angular domain and thus reduces the number of RF chains due to the angular sparsity of the signals. Since it does not require any phase shifters, the total cost is reduced compared with hybrid beamforming architecture. The other approach is based on digital beamforming which involves *multi-layer signal processing* [7]. Although the number of RF chains is not reduced, it can leverage the different levels of CSI and hardware to reduce the processing complexity and pilot overhead.

In essence, all the aforementioned solutions aim at providing comparable performance as full digital preprocessing systems with limited number of RF chains or reduced complexity in massive MIMO systems. See Fig. 1 for an illustration. The current literature mainly focuses on the right side of the figure, i.e., the antenna to beam mapping and beam selection schemes, which leverages the angular domain power sparsity of the channel to transform the signals from the antenna domain to the beam domain and hence maps to RF chains. On the other hand, the left side of the Fig. 1 which represents user-beam mapping, is scantily treated in the literature. The user-beam mapping essentially deals with **user scheduling in the beam domain**. Unlike the previous extensive user scheduling related work, e.g., in [8], [9], the user scheduling in the beam domain has its uniqueness, that is the user scheduling problems are tangled with the beam selection. In reality, due to the angular sparsity of the massive MIMO

channel [10], the beams, which represent the signal directions, are strongly related to the users, in the sense that each beam usually contains signals of very few (possibly one) users. Therefore, the user scheduling and beam selection have to be jointly considered. Furthermore, it is proposed that beam-based downlink scheduling should be performed *only* based on the **statistical CSI**, e.g., second-order CSI, whereas the existing literature often assumes instantaneous CSI is available [11]. The reason of only assuming statistical CSI is two-fold. First, the statistical CSI is much less costly to obtain than instantaneous CSI in terms of overhead, attributing to the fact that statistical CSI changes much more slowly. Moreover, if we consider the C-RAN system, the beamforming module is integrated with the remote radio heads (RRHs) and hence limited computation power is expected which prevents us from using complicated channel estimation schemes.

In this paper, we aim to address the user and beam joint scheduling problem in beam-based massive MIMO downlinks. The contributions are as follow. We formulate the problem based on the Lyapunov-drift optimization framework. An optimal scheduling policy is proposed thereby to achieve optimum utilities. The optimality proof is given which shows the achieved utility is arbitrarily close to the optimum. To address the queue weighted sum rate maximization problem arisen in optimizing the Lyapunov-drift, which is a mixed integer programming (MIP) problem, the block-coordinated-update-based (BCU-based) algorithm which deals with the continuous convex relaxation of the MIP problem is proposed. An iterative water-filling based approach is also proposed to reduce the number of iterations. In order to implement the algorithm based on the statistical CSI, a deterministic equivalent of the capacity of the downlink broadcast channel in the large antenna array regime is derived.

## II. SYSTEM MODEL AND PROBLEM FORMULATIONS

### A. Signal Model

The single-cell system downlink is considered in this paper, where a BS with  $M$  co-located antennas transmits to  $N_t$  single-antenna users. The BS has  $K$  RF chains ( $K \leq M$ ). Assuming narrow-band and time-invariant channels, the receive signal of user- $n$  is,

$$y_n = \mathbf{h}_n^\dagger \mathbf{x} + n_n, \quad (1)$$

where  $\mathbf{h}_n$  is an  $M$ -dimensional channel vector,  $\mathbf{x}$  is the downlink transmit signals, and  $n_n$  denotes the i.i.d. Gaussian additive noise with unit variances. The downlink channel matrix is denoted by  $\mathbf{H} = [\mathbf{h}_1, \mathbf{h}_2, \dots, \mathbf{h}_{N_t}]^\dagger$ . The transmit signal including beamforming can be written as  $\mathbf{x} = \mathbf{B}_a \mathbf{B}_d \mathbf{s}$ , where  $\mathbf{s}$  denotes the  $N_s$ -dimensional data symbols for the scheduled spatial-multiplexing users ( $N_s \leq K$  for complete inter-user-interference elimination, and obviously  $N_s \leq N_t$ ), and hence the number of scheduled users is also  $N_s$ . The RF (analog) beamforming, which can be realized by the lens antenna array and beam selection in beamspace MIMO or general analog beamforming in hybrid beamforming architectures, is denoted by  $\mathbf{B}_a$  with dimension  $N_t \times K$ . The digital precoding at

baseband is denoted by the matrix  $\mathbf{B}_d$  with dimension  $K \times N_s$ . On account of the analog beamforming, the effective channel observed from the digital baseband is  $\bar{\mathbf{H}} = \mathbf{H} \mathbf{B}_a$ , where the effective channel vector corresponding to user- $n$  is denoted by  $\bar{\mathbf{h}}_n$  and  $\bar{\mathbf{H}} = [\bar{\mathbf{h}}_1, \dots, \bar{\mathbf{h}}_{N_t}]^\dagger$ .

The RF beamforming considered in the paper is the widely adopted directional beamforming scheme, which is in line with the beamspace MIMO architecture, and hence  $\mathbf{B}_a = \mathbf{B}_{\text{DFT}} \boldsymbol{\Sigma}_b$ , where  $\mathbf{B}_{\text{DFT}}$  is the equivalent discrete-Fourier-transform (DFT) matrix (or Kronecker product of DFT matrices for two-dimensional antenna array). The beam selection decision is denoted by the diagonal matrix  $\boldsymbol{\Sigma}_b$  whose entries are binary, i.e.,  $(\boldsymbol{\Sigma}_b)_{i,i} \in \{0, 1\}$ ,  $\forall i$ .

### B. Problem Formulations

The long-time average rate of user  $n$  is denoted by  $\bar{R}_n = \lim_{T \rightarrow \infty} \frac{1}{T} \sum_{t=0}^{T-1} R_n(\mathbf{H}(t), \pi(t))$ , where  $R_n(\mathbf{H}(t), \pi(t))$  denotes the instantaneous rate of user  $n$  at time  $t$ , given the channel coefficients  $\mathbf{H}(t)$  and control (user scheduling and beam selection as far as the paper is concerned)  $\pi(t)$ . Based on ergodicity,  $\bar{R}_n = \mathbb{E}\{R_n(\mathbf{H}, \pi)\}$ ,  $\forall n$ , where the expectation is taken over channel coefficients  $\mathbf{H}(t)$  and possibly  $\pi(t)$  when a stochastic control policy is considered. The achievable ergodic rate region can be characterized as  $\mathcal{R} = \text{coh} \bigcup_{\pi \in \mathcal{X}} \{\bar{\mathbf{R}} : 0 \leq \bar{R}_n \leq \mathbb{E}[R_n(\mathbf{H}, \pi)]\}$ , where  $\bar{\mathbf{R}}$  is a  $N_t$ -dimensional region,  $\bar{R}_n$  is its  $n$ -th component, and ‘‘coh’’ denotes the closure of a convex hull. The set of all feasible scheduling policies is denoted by  $\mathcal{X}$ . The problem is formulated as

$$\mathbf{P1}: \quad \underset{\boldsymbol{\Sigma}_u, \boldsymbol{\Sigma}_b}{\text{maximize}} U(\bar{\mathbf{R}}), \text{ s.t., } \bar{\mathbf{R}} \in \mathcal{R}, \quad (2)$$

where  $\boldsymbol{\Sigma}_u$  is a diagonal matrix denoting user scheduling decision at time  $t$ , i.e.,  $s_i \triangleq (\boldsymbol{\Sigma}_u)_{i,i} \in \{0, 1\}$ , and  $b_i \triangleq (\boldsymbol{\Sigma}_b)_{i,i}$  denotes the beam selection decision. The network utility function  $U(\bar{\mathbf{R}})$  is defined as a function of the long-time average rate for each user, e.g.,

$$U_{\text{SUM}}(\bar{\mathbf{R}}) = \sum_n \bar{R}_n \quad (3)$$

for sum rate maximization,

$$U_{\text{PFS}}(\bar{\mathbf{R}}) = \sum_n \log(\bar{R}_n + c_n) \quad (4)$$

for proportional-fairness scheduling (PFS) [12], where  $c_n$ 's are non-negative constants to regularize the logarithm expressions, and typical value is  $c_n = 0$ ,  $\forall n$  for exact PFS or  $c_n = 1$ ,  $\forall n$  to ensure positive objective function value which is a mathematical convenience. Basic properties of the utility function  $U(\bar{\mathbf{R}})$  are required, e.g., concavity and non-decreasing [13].

## III. OPTIMAL BEAM-BASED DOWNLINK SCHEDULING POLICY

To address the optimization in (2), it is found that two severe challenges exist. First, the ergodic capacity region, i.e.,  $\mathcal{R}$ , does not yield a closed-form expression. The work in

[14] [15] characterizes the broadcast channel (BC) capacity region and the duality between BC and multiple-access-channel (MAC) in the sense of both capacity region and outage probability. However, no closed-form expressions are presented except for capacity bounds [16]. Secondly, the scheduling and beam selection decisions should be made *dynamically* to match the channel variations and user traffic in time. To address these issues, we seek to leverage a powerful tool of Lyapunov-drift optimization which is shown to have superior performance compared to static solutions.

#### A. Lyapunov-Drift based Network Utility Maximization

To maximize the network utility function in (2), the transmission need of each user, which is determined by the transmission history and utility function, is represented by virtual queues. The arrival process is designed to reflect the transmission need and a max-weight algorithm is applied to stabilize the system. Specifically, let  $Q_n(t)$  denote the virtual queue length in bits of user  $n$  at the beginning of  $t$ -th scheduling step, let  $a_n(t)$  denote the number of arrival bits which are designed later, and let  $\mu_n(t)$  denote the allocated number of service bits to queue- $n$ , which equals the allocated number of service bits between scheduling steps. The queuing dynamics are written as  $Q_n(t+1) = Q_n(t) - \tilde{\mu}_n(t) + a_n(t)$ , where  $\mu_n(t) = \sum_{\tau=1}^T R_n(\mathbf{H}(t+\tau), \pi(t+\tau))$ , and  $\tilde{\mu}_n(t) = \min\{Q_n(t), \mu_n(t)\}$  denotes the number of actual service bits, considering the circumstances that sometimes the queue is emptied given the amount of allocated service bits. The optimal beam-based downlink scheduling policy at a given scheduling time  $t$ , i.e., a dynamic (possibly randomized) policy which achieves the solution to (2), can be described as below.

**Admission control:** For virtual queue  $\mathbf{Q}(t) = [Q_1(t), \dots, Q_{N_t}(t)]$ , and  $q_n = Q_n(t)$  denotes the virtual queue state at the scheduling time. Let the arrival  $\mathbf{a}$  be the solution of

$$\mathbf{P2:} \quad \underset{\mathbf{a}}{\text{maximize}} \quad VU(\mathbf{a}) - \mathbf{a}(t)^T \mathbf{Q}(t), \quad (5)$$

$$\text{s.t.}, \quad 0 \leq a_n \leq A_{\max}, \quad \forall n, \quad (6)$$

where  $V$  and  $A_{\max}$  are pre-defined constants<sup>2</sup>.

**Scheduling:** Given the arrival process determined above, the service, i.e., the scheduling, is based on the solution of the following problem:

$$\mathbf{P3:} \quad \underset{\Sigma_u, \Sigma_b, \mathbf{p}}{\text{maximize}} \quad \sum_{n=1}^{N_t} \left[ q_n s_n \sum_{\tau=1}^T R_n(\mathbf{H}(\tau), \Sigma_u, \Sigma_b, \mathbf{p}) \right] \quad (7)$$

$$\text{s.t.}, \quad \sum_{n=1}^{N_t} p_n \leq P, \quad \sum_{i=1}^{N_t} s_i = N_s, \quad \sum_{i=1}^M b_i = K \quad (8)$$

where  $s_i \in \{0, 1\}$ ,  $b_i \in \{0, 1\}$ ,  $\mathbf{p} = [p_1, \dots, p_{N_t}]$  denotes the transmit power of  $N_t$  user data streams and hence  $P$  in (8) is the sum power constraint. The scheduling decisions are

denoted by binary variables  $s_i$  and  $b_i$ .  $\Sigma_b$  and  $\Sigma_u$  are diagonal matrices consisting of  $s_i$  and  $b_i$ . The downlink instantaneous transmission rate  $R_n(\mathbf{H}(t), \Sigma_u, \Sigma_b, \mathbf{p})$  is a function of the downlink transmit power allocation, channel coefficients, and scheduling decisions. It is observed that the admission control problem **P2** is a convex problem and hence relatively easier to solve, e.g., for PFS, the optimum admission control is given by  $a_n = \min\left\{\frac{V}{q_n}, A_{\max}\right\}$ . However, the problem **P3** is an MIP problem, which is NP-complete [17]. Therefore, we seek to relax the binary constraints in the next subsection. Before diving into details on solving **P3**, we assume the optimum solutions to both problems are obtained for the moment, which is denoted by  $\pi^*$ . The optimality of the scheduling policy is established in the following theorem.

*Theorem 1:* Denote  $\bar{\mathbf{R}}^* = \arg \max_{\bar{\mathbf{R}} \in \mathcal{R}} U(\bar{\mathbf{R}})$ . Suppose the transmission rate is bounded, i.e.,  $0 \leq \bar{R}_n \leq R_{\max}, \forall n$ , the utility function  $U(\cdot)$  is concave and entry-wise non-decreasing which is bounded on  $[0, R_{\max}]$ , and that the channel coefficients  $\mathbf{H}(t)$  are i.i.d. over different scheduling period, then based on the scheduling algorithm resulting from **P2** and **P3**, the following conditions are met.

$$\liminf_{\tau \rightarrow \infty} U\left(\frac{1}{\tau} \sum_{t=0}^{\tau-1} \mathbb{E}[\mathbf{R}(t)]\right) \geq U(\bar{\mathbf{R}}^*) - C/V, \quad (9)$$

$$\lim_{\tau \rightarrow \infty} \frac{\mathbb{E}[Q_n(\tau)]}{\tau} = 0, \quad \forall n \quad (10)$$

i.e., the utility function of the time-averaged transmission rate based on the scheduling given above is within a constant (arbitrary small if  $V$  is large) to the optimum and the virtual queues are mean-rate-stable, where  $C$  is a constant related to  $A_{\max}$  which, along with  $V$ , is defined in **P2**.

*Proof:* The proof is based on the Lyapunov-drift framework. The details are omitted due to lack of space. ■

#### IV. BLOCK COORDINATE UPDATE BASED JOINT BEAM AND USER SCHEDULING

This section is dedicated to solving the scheduling problem of **P3** only based on the knowledge of statistical CSI. The last section establishes the optimality of the proposed beam-based scheduling algorithm given the solutions of **P2** (generally easy to solve) and **P3**. However, due to the NP-hardness of **P3**, explicit solutions are hard to attain. More importantly, it is proposed that the scheduling decisions of **P3** should only rely on statistical CSI, rendering the solution even more intractable. Towards this end, an algorithm based on solving the convex relaxation of the original problem leveraging the BCU technique and random matrix theory is proposed. First, **P3** is first transformed for better exposition based on the uplink-downlink duality [18]. Observing the objective function in **P3**, the transmission rate is evaluated by the MIMO broadcast channel capacity. The following Proposition 1 derives an implicit expression of the objective function and gives an asymptotic result such that the optimization is only dependent on statistical CSI which in this case is the channel correlation matrices.

<sup>2</sup>For Typical values,  $V$  and  $A_{\max}$  can be approximately 100-fold of the service rate.

*Proposition 1:* In the large system regime, i.e.,  $K \rightarrow \infty$  and  $K/N_s \rightarrow \beta$ , the queue-weighted downlink achievable rate in **P3** is asymptotically equivalent to

$$\begin{aligned} & \mathcal{D}(\mathbf{Q}, \mathbf{R}_1, \dots, \mathbf{R}_{N_t}, \Sigma_u, \Sigma_b, \mathbf{p}) \\ \triangleq & \sum_{n=1}^{N_t} q_n T \log \left( 1 + p_n s_n \text{tr} \left[ \Sigma_b \mathbf{B}_{\text{DFT}}^\dagger \mathbf{R}_n \mathbf{B}_{\text{DFT}} \Sigma_b \right. \right. \\ & \left. \left. \left( \frac{1}{M} \sum_{j=1}^{n-1} \frac{p_j s_j \Sigma_b \mathbf{B}_{\text{DFT}}^\dagger \mathbf{R}_j \mathbf{B}_{\text{DFT}} \Sigma_b}{1 + e_{n,j}} + \mathbf{I} \right)^{-1} \right] \right), \end{aligned} \quad (11)$$

where  $q_i$ 's are arranged in non-increasing order, and  $e_{n,i}$  is the unique solution of the following equations.

$$\begin{aligned} e_{n,i} = & \text{tr} \left[ \Sigma_b \mathbf{B}_{\text{DFT}}^\dagger \mathbf{R}_i \mathbf{B}_{\text{DFT}} \Sigma_b \right. \\ & \left. \left( \frac{1}{M} \sum_{j=1}^{n-1} \frac{p_j s_j \Sigma_b \mathbf{B}_{\text{DFT}}^\dagger \mathbf{R}_j \mathbf{B}_{\text{DFT}} \Sigma_b}{1 + e_{n,j}} + \mathbf{I} \right)^{-1} \right] \end{aligned} \quad (12)$$

*Proof:* The proof is based on random matrix theory, and the details are omitted due to lack of space. ■

#### A. Convex Relaxation

Although the original MIP **P3** which is NP-complete, it can be transformed to a multi-convex problem based on Proposition 1 by relaxing the binary constraints to real-value constraints. The convex relaxation of an MIP is a widely-used technique to achieve near-optimal solution to the original problem [19] [20]. The relaxed version of **P3** is stated below.

$$\begin{aligned} \mathbf{P4}: & \text{maximize}_{\Sigma_b, \mathbf{w}} \mathcal{D}(\mathbf{Q}, \mathbf{R}_1, \dots, \mathbf{R}_{N_t}, \mathbf{I}_{N_t}, \Sigma_b, \mathbf{w}) \\ & \text{s.t.}, \sum_{n=1}^{N_t} w_n \leq P, \sum_{i=1}^M b_i = K, 0 \leq b_i \leq 1, \forall i, \end{aligned} \quad (13)$$

where  $w_n = p_n s_n$ , and  $\mathcal{D}$  is defined in (11). Without loss of generality, define the optimum solution of **P4** as  $b_i^*$ 's and  $w_n^*$ 's, respectively. Then the scheduling decision is to schedule the beams and users corresponding to the largest  $K$   $b_i^*$ 's and  $N_s$   $w_n^*$ 's, respectively. It is observed that **P4** is a multi-convex problem since the objective function is concave in both  $\Sigma_b$  and  $\mathbf{w}$  [20]. In view of this, the following Algorithm 1, which bases upon the BCU technique is proposed.

The basic idea of the proposed BCU-based user and beam joint scheduling is that an iterative method which cyclically optimizes user scheduling and beam selection with the other fixed is guaranteed to converge to the global optimum of **P4**. In order to accelerate the iteration, an iterative water filling approach is adopted.

*Convergence of the BCU-based algorithm:* The convergence to the global optimum is due to the convergence results of the BCU algorithm [21]. The details of the proof is omitted due to lack of space.

---

#### Algorithm 1: BCU-Based Scheduling

---

- 1 **Initialization:**  $\Sigma_b^{(0)} = \mathbf{I}_M$ ;
  - 2 **Iteration:** for  $t = 1 : T$  do
  - 3     **User scheduling update based on iterative water filling:**  $\forall n \in [1, N_t], \omega_n^{(0)} = P/N_t$ ; for  $t_w = 1 : T_w$  do
  - 4         Compute for each  $n$ ,  
 $\beta_n^{(t_w)} = \text{tr} \left[ \Sigma_b^{(t-1)} \mathbf{B}_{\text{DFT}}^\dagger \mathbf{R}_n \mathbf{B}_{\text{DFT}} \Sigma_b^{(t-1)} \right.$   
 $\left. \left( \frac{1}{M} \sum_{j=1}^{n-1} \frac{w_j^{(t_w-1)} \mathbf{R}_j}{1 + e_{n,j}} + \mathbf{I} \right)^{-1} \right]$ , where  $e_{n,i}$  is  
the unique solution of the equations in (12);
  - 5         Apply the classical water filling algorithm with  
water levels defined by  $\beta^{(t_w)}$   
 $\gamma^{(t_w)} = \arg \max_{\sum_n \gamma_n \leq P, \gamma \geq 0} \sum_{n=1}^{N_t} q_n \log \left( 1 + \gamma_n \beta_n^{(t_w)} \right)$
  - 6         Update  $\omega$  as  
 $\omega^{(t_w)} = (1 - 1/M) \omega^{(t_w-1)} + (1/M) \gamma^{(t_w)}$
  - 7         **if**  $\|\omega_n^{(t_w)} - \omega_n^{(t_w-1)}\| < \epsilon$  **then**
  - 8              $\mathbf{w}^{(t)} = \omega^{(t_w)}$ , break;
  - 9         **Beam selection:** Solve for  $\Sigma_b^{(t)}$ , which is the  
solution to the convex optimization problem of **P4**  
with  $\mathbf{w} = \mathbf{w}^{(t)}$ .
  - 10         **Stopping criterion:** **if**  $\|\mathbf{w}^{(t)} - \mathbf{w}^{(t-1)}\| < \epsilon_1$  **and**  
 $\|\Sigma_b^{(t)} - \Sigma_b^{(t-1)}\| < \epsilon_2$  **then**
  - 11              $\mathbf{w}_{\text{opt}} = \mathbf{w}^{(t)}$ ,  $\Sigma_{b,\text{opt}} = \Sigma_b^{(t)}$ , break;
  - 12 **Output:** The scheduling user set is the users with the  
largest  $N_t$  values in  $\mathbf{w}_{\text{opt}}$ . The selected beams are the  
ones with the largest  $K$  values in the diagonal entities  
in  $\Sigma_{b,\text{opt}}$ .
- 

*Remark 1:* The Algorithm 1 can solve the joint user scheduling and beam selection problem based on only the statistical CSI. Therefore, it is applicable before the channel estimations. After the system selects the users and beams, the instantaneous channel estimations can be implemented and hence digital precoding and decoding can follow. This is in line with the multi-layer signal processing concept proposed in, e.g., [3], which proposes that the pre-beamforming should be done based on only channel statistics to save RF chains, complexity and system overhead. Additionally, the optimality of the BCU-based algorithm regarding the convex relaxation problem is proved.

#### V. SIMULATION RESULTS

In this section, the simulation results are presented to validate the performance gain of the proposed BCU-based scheduling scheme. The number of MPCs for each user is  $L_n = 3, \forall n$  (including one LoS MPC). The ULA is used in the simulations and the amplitude of the LoS MPC is 10 times the one of the NLoS MPCs. The DoAs of the signals

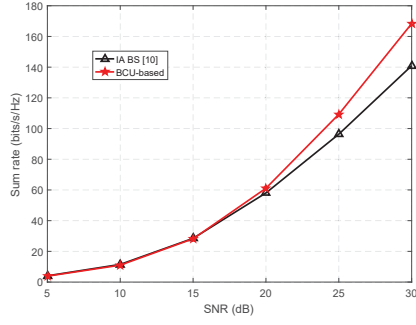


Fig. 2. Sum rate comparisons with identical user pathloss. The number of BS antennas is 256, the number of users is 100, the number of scheduled users and beams are both 40.

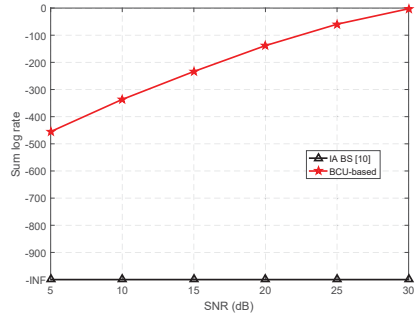


Fig. 3. Sum log rate comparisons with non-identical user pathloss. The users with non-identical pathloss are generated with distances i.i.d. uniformly distributed between 20 m to 200 m. The number of BS antennas is 256, the number of users is 100, the number of scheduled users and beams are both 40.

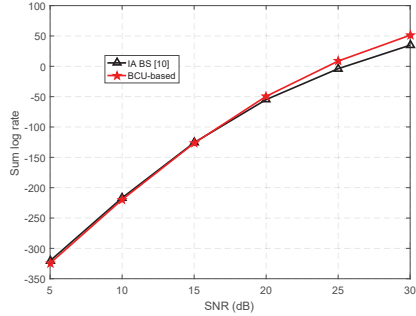


Fig. 4. Sum log rate comparisons with identical user pathloss. The users with non-identical pathloss are generated with distances i.i.d. uniformly distributed between 20 m to 200 m. The number of BS antennas is 256, the number of users is 100, the number of scheduled users and beams are both 40.

are generated from i.i.d. uniform distributions. The antenna spacing  $d = \lambda/2$ , where  $\lambda$  denotes the carrier wavelength. In some of the following cases where users' pathlosses are not identical, the distances of users are generated based on an i.i.d. uniform distributions from 30 to 200 meters and the pathloss  $\gamma_n$  is  $\gamma_n = \left(\frac{d_n}{d_0}\right)^{-\gamma}$ , where  $\gamma = 2$  which is in line with millimeter-wave channel measurements [22] and  $d_0$  is some reference point distance. The regularized zero-forcing (RZF) precoder is adopted for system evaluation, i.e., define  $\mathbf{K}_{\text{RZF}} = (\bar{\mathbf{H}}^\dagger \bar{\mathbf{H}} + M\alpha \mathbf{I}_M)^{-1}$ . The RZF precoding matrix is

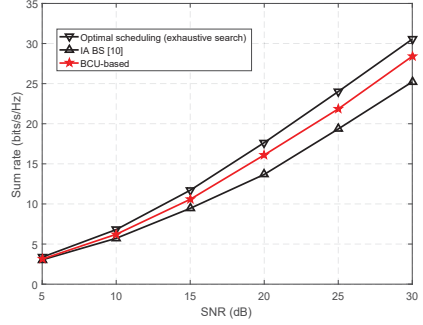


Fig. 5. Sum rate comparisons with optimal scheduling. The number of BS antennas is 8, the number of users is 8, the number of scheduled users and beams are both 4.  $c_n = 1, \forall n$  as in (4).

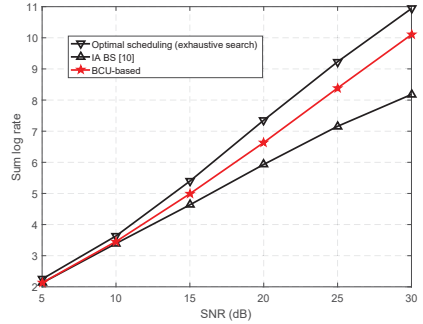


Fig. 6. Sum log rate comparisons with optimal scheduling. The number of BS antennas is 8, the number of users is 8, the number of scheduled users and beams are both 4.  $c_n = 1, \forall n$  as in (4).

expressed as  $\mathbf{B}_d = \zeta \mathbf{K}_{\text{RZF}} \bar{\mathbf{H}}^\dagger$ , where  $\zeta$  is a normalization scalar to fulfill the power constraint, and  $\alpha$  is the regularization factor. Although RZF precoder is not the optimal coding scheme for Gaussian broadcast channel (dirty-paper-coding with minimum-mean-square-error (MMSE) precoder is proved for optimality but the usage is limited in reality due to high complexity), it can achieve full degree-of-freedom (DoF) in the high signal-to-noise-ratio (SNR) region and easy to implement. In the simulations,  $\alpha = N_s/\rho$ , where  $\rho$  is the SNR. The user instantaneous rate is calculated by the Shannon formula. The block fading model is adopted, where the channel stay constant for 10 time slots and evolves to another constant based on an i.i.d. distribution, which in this case involves the generation of the random MPC phase and amplitude. The simulation runs for 10000 such blocks and calculate the time-averaged downlink transmission rates.

The constants used in the Lyapunov-drift optimization are set to be  $V = A_{\text{max}} = 10^2 r_{\text{max}}$ , where  $r_{\text{max}}$  is the maximum rate of the users. The  $\epsilon$ ,  $\epsilon_1$  and  $\epsilon_2$  in the stopping criterion in BCU-based algorithm are set to be  $10^{-2} \rho/K$ . In comparisons, the IA-based scheduling (IA BS) in [11] is also simulated, which is based on the instantaneous CSI, whereas the our proposed schemes are based on statistical CSI.

In Fig. 2, our proposed BCU-based algorithm as in Algorithm 1 is compared with the IA BS algorithm. Since the sum rates are considered, the utility function in (3) is adopted.

The scheduled user set in IA BS is a greedy choice for comparisons, which is the users with the most channel power. It is observed that by jointly considering user scheduling and beam selection, the sum rate performance can be improved in the high-SNR regime. The reason is that in the high-SNR regime, the interference is dominating the performance, and thus by jointly considering the user scheduling and beam selection by the proposed schemes, the interference is better suppressed.

Considering the other utility functions, e.g., the PFS utility function in (4), the performance advantage is more obvious as shown in Fig. 3. In the IA BS, the beams with stronger channels are always preferred, corresponding to users with small pathloss. Therefore, the fairness among users is neglected. In the proposed joint scheduling schemes, the Lyapunov admission control in **P2** utilizes a virtual queue to control the fairness among users. Note that since there are always unscheduled users in the IA BS due to their small pathloss, the sum log rate of the IA BS scheme in this case is negative infinity. Alternatively, it can be interpreted as the geometrical mean of the user rates is zero. In Fig. 4, the performance with identical pathloss is presented. Due to the fact that users have identical large-scale fading, and hence equal probability to have good channels, the performance advantage over IA BS is relatively small in this case.

In order to show the performance loss of the proposed scheme compared with optimal scheduling, an exhaustive search over all the feasible user and beam sets is conducted to solve **P3** and the optimal scheduling performance is obtained accordingly. Due to the prohibitive high complexity of exhaustive search, we consider a small-scale problem where there are 8 users, 8 BS antennas and 4 scheduled users and beams. The results are shown in Fig. 5 and 6, which show sum-rate and sum-log-rate optimizations, respectively. It is observed that in general the proposed schemes can achieve near-optimal performance. The sum log rate is plotted with  $c_n = 1, \forall n$  in (4) in this case to avoid negative infinity value for the IA BS for comparisons.

## VI. CONCLUSIONS

In this paper, the BCU-based scheduling scheme is proposed to address the joint user and beam scheduling problem in beam-based massive MIMO systems based only on statistical CSI. The problem is formulated under the Lyapunov-drift optimization framework. In order to solve the weighted rate maximization problem therein, which is an NP-hard problem, the BCU-based scheme leverages the convex relaxation of the problem and adopts an iterative water-filling approach. It is proved that the BCU-based scheduling scheme iteratively converges to optimum of the relaxed problem. In the simulations, it is shown that the proposed scheme can achieve near-optimal performance and outperform the state-of-the-art beam selection schemes, with utilities such as sum rate and proportional fairness. It is found that when the user scheduling and beam selection are jointly optimized, especially when user

fairness is considered, the performance can be significantly improved.

## REFERENCES

- [1] T. Marzetta, "Noncooperative cellular wireless with unlimited numbers of base station antennas," *IEEE Trans. Wireless Commun.*, vol. 9, pp. 3590–3600, Nov 2010.
- [2] E. G. Larsson, O. Edfors, F. Tufvesson, and T. L. Marzetta, "Massive MIMO for next generation wireless systems," *IEEE Commun. Mag.*, vol. 52, pp. 186–195, Feb. 2014.
- [3] R. W. Heath, N. Gonzalez-Prelcic, S. Rangan, W. Roh, and A. M. Sayeed, "An overview of signal processing techniques for millimeter wave MIMO systems," *IEEE J. Sel. Top. Signal Process.*, vol. 10, pp. 436–453, Apr. 2016.
- [4] Z. Jiang, S. Zhou, R. Deng, Z. Niu, and S. Cao, "Pilot-data superposition for beam-based FDD massive MIMO downlinks," *IEEE Commun. Letters*, vol. 21, pp. 1357–1360, Jun 2017.
- [5] A. F. Molisch, V. V. Ratnam, S. Han, Z. Li, S. L. H. Nguyen, L. Li, and K. Haneda, "Hybrid beamforming for massive MIMO-A survey," *arXiv preprint arXiv:1609.05078*, 2016.
- [6] J. Brady, N. Behdad, and A. M. Sayeed, "Beamspace MIMO for millimeter-wave communications: System architecture, modeling, analysis, and measurements," *IEEE Trans. Antennas and Propagat.*, vol. 61, pp. 3814–3827, Jul. 2013.
- [7] Z. Jiang, A. Molisch, G. Caire, and Z. Niu, "Achievable rates of FDD massive MIMO systems with spatial channel correlation," *IEEE Trans. Wireless Commun.*, vol. 14, pp. 2868–2882, May 2015.
- [8] T. Yoo and A. Goldsmith, "On the optimality of multi-antenna broadcast scheduling using zero-forcing beamforming," *IEEE J. Select. Areas Commun.*, vol. 24, pp. 528–541, Mar. 2006.
- [9] H. Shirani-Mehr, G. Caire, and M. J. Neely, "MIMO downlink scheduling with non-perfect channel state knowledge," *IEEE Trans. Commun.*, vol. 58, pp. 2055–2066, July 2010.
- [10] C. U. Bas, R. Wang, D. Psychoudakis, T. Henige, R. Monroe, J. Park, J. Zhang, and A. F. Molisch, "A real-time millimeter-wave phased array MIMO channel sounder," *arXiv preprint arXiv:1703.05271*, 2017.
- [11] X. Gao, L. Dai, Z. Chen, Z. Wang, and Z. Zhang, "Near-optimal beam selection for beamspace mmwave massive MIMO systems," *IEEE Commun. Letters*, vol. 20, pp. 1054–1057, May 2016.
- [12] H. J. Kushner and P. A. Whiting, "Convergence of proportional-fair sharing algorithms under general conditions," *IEEE Trans. Wireless Commun.*, vol. 3, pp. 1250–1259, Jul 2004.
- [13] M. J. Neely, E. Modiano, and C. P. Li, "Fairness and optimal stochastic control for heterogeneous networks," *IEEE/ACM Trans. Networking*, vol. 16, pp. 396–409, Apr 2008.
- [14] H. Weingarten, Y. Steinberg, and S. Shamai, "The capacity region of the Gaussian multiple-input multiple-output broadcast channel," *IEEE Trans. Inform. Theory*, vol. 52, pp. 3936–3964, Sep. 2006.
- [15] W. Yu, "Uplink-downlink duality via minimax duality," *IEEE Trans. Inform. Theory*, vol. 52, pp. 361–374, Feb. 2006.
- [16] Z. Jiang, S. Zhou, and Z. Niu, "Capacity bounds of downlink network MIMO systems with inter-cluster interference," in *IEEE Global Communications Conference (GLOBECOM)*, pp. 4612–4617, Dec. 2012.
- [17] R. M. Karp, "Reducibility among combinatorial problems," in *Complexity of computer computations*, pp. 85–103, Springer, 1972.
- [18] W. Yu and J. Cioffi, "Sum capacity of Gaussian vector broadcast channels," *IEEE Trans. Inform. Theory*, vol. 50, pp. 1875–1892, Sep. 2004.
- [19] S. E. Nai, W. Ser, Z. L. Yu, and H. Chen, "Beampattern synthesis for linear and planar arrays with antenna selection by convex optimization," *IEEE Trans. Antennas and Propag.*, vol. 58, pp. 3923–3930, Dec 2010.
- [20] A. Dua, K. Medepalli, and A. J. Paulraj, "Receive antenna selection in MIMO systems using convex optimization," *IEEE Trans. Wireless Commun.*, vol. 5, pp. 2353–2357, Sep. 2006.
- [21] P. Tseng, "Convergence of a block coordinate descent method for nondifferentiable minimization," *Journal of Optimization Theory and Applications*, vol. 109, no. 3, pp. 475–494, 2001.
- [22] T. S. Rappaport, S. Sun, R. Mayzus, H. Zhao, Y. Azar, K. Wang, G. N. Wong, J. K. Schulz, M. Samimi, and F. Gutierrez, "Millimeter wave mobile communications for 5G cellular: It will work!," *IEEE Access*, vol. 1, pp. 335–349, 2013.